

Express Mailing Label No. EL 733 633 177 US

PATENT APPLICATION
Docket No. 14113.3.2.2

UNITED STATES PATENT APPLICATION

of

Michael R. Ohran

for

Mirroring Network Data to Establish

Virtual Storage Area Network

WORKMAN, NYDEGGER & SEELEY
A PROFESSIONAL CORPORATION
ATTORNEYS AT LAW
1000 EAGLE GATE TOWER
60 EAST SOUTH TEMPLE
SALT LAKE CITY, UTAH 84111

BACKGROUND OF THE INVENTION

1. Related Applications

This application is a continuation-in-part of U.S. Patent Application Serial No. 09/271,585, entitled "Operation of Standby Server to Preserve Data Stored By a Network Server," filed March 18, 1999, which is a continuation of U.S. Patent Application Serial No. 08/848,139, filed April 28, 1997, entitled "Method for Rapid Recovery from a Network File Server Failure Including Method for Operating Co-Standby Servers," now issued as U.S. Patent No. 5,978,565. The foregoing patent and patent application are incorporated herein by reference.

2. The Field of the Invention

This invention relates to network server computer systems, and in particular an improvement to the methods used to recover from a computer failure in a system that provides a virtual storage area network, in which multiple server computers access the same network data.

3. Background and Related Art

In a network server computer system, there are a plurality of personal computers or user workstations that are usually supported by two or more servers. In order to provide continuous operation of these computer systems, it is necessary for the computer system to provide a method for overcoming faults and failures that often occur within the network server computer system. This is generally done by having redundant computers and mass storage devices, such that a backup server computer or disk drive is immediately available to take over in the event of a fault or failure of a primary server computer or disk drive.

1 A technique for implementing a fault-tolerant computer system is described in
2 Major et al., United States Patent No. 5,157,663. In particular, Major provides a redundant
3 network file server system capable of recovering from the failure of either the computer or
4 the mass storage device of one of the file servers. The file server operating system is run
5 on each computer system in the network file server, with each computer system
6 cooperating to produce the redundant network file server. This technique has been used by
7 Novell, of Provo, UT, to implement its SFT-III fault-tolerant file server product.

8 More recently, fault-tolerant networks known as "storage area networks" have been
9 developed. A storage area network ("SAN") connects multiple servers of an enterprise
10 network with a common or shared storage node to store and access network data. In the
11 case of a failure of one of the servers, the other servers can perform network services that
12 would otherwise have been provided by the failed server.

13 Figure 1 illustrates a typical architecture of a network system that includes a
14 conventional storage area network. Figure 1 illustrates three server computers 110, 120,
15 and 130 that provide network services for network 101. Although three servers are
16 illustrated in Figure 1, network 101 may include as few as two servers or more servers than
17 are shown in Figure 1. This variable number of server computers depends upon the
18 individual needs of the network being served. For example, a large organization may
19 require the use of several server computers, likewise a smaller organization might simply
20 require two server computers.

21 In this configuration, user workstations (or personal computers) 102 are connected
22 to network 101 and have access to server computers 110, 120, and 130. Each user
23 workstation is generally associated with a particular sever computer, although, in a
24 network system that includes a storage area network, any server can provide substantially

any network services for any workstation, as needed. A user, at a user workstation 102, issues requests for operations, such as read, write, etc., which are transmitted to the associated server computer, 110, 120, or 130, which then performs the requested operation using I/O drivers 113, 123, and 133. Servers 110, 120, and 130 perform data operations on network data that is stored in disks 142 of shared storage node 140. Each server 110, 120, and 130 has access to any network data stored at shared storage node 140, subject to policing protocol described below. The storage area network of Figure 1 includes the physical communication infrastructure and the protocols that enable server computers 110, 120, and 130 to operate with shared storage node 140.

Each server computer includes software representing a policing protocol module 111, 121, 131, that cooperates with the policing protocol modules of the other server computers to implement a policing protocol. The policing protocol prevents data corruption by controlling the performance of requested operations. For example, the policing protocol implemented by modules 111, 121, and 131 may allow a server to respond to read operation requests at any time, but may permit only one server computer at a time to perform a write operation request.

One advantage of SANs is that all server computers have access to all network data through the shared storage node. If one server experiences a failure, workstations can bypass the failed server and issue operation requests to other servers. The shared storage node prevents the need for mirroring data between multiple storage nodes associated with different servers. However, storage area networks have at least two significant liabilities that have prevented them from becoming fully accepted in the marketplace and make them unsuitable for many customers.

24

1 First, SANs require specialized hardware, namely, the shared storage node. Many
2 potential users of storage area networks find the cost of purchasing and maintaining a
3 shared storage node prohibitive. In practice, many users of SANs are large corporations or
4 other enterprises that have relatively large networks with large numbers of servers.
5 Enterprises that have the need for only two or three servers may not find it cost-effective to
6 implement a storage area network.

7 Second, although SANs are tolerant of failures of network servers, they are not well
8 suited for responding or protecting against other hardware failures. For example, because
9 a storage area network uses a single shared storage node, any failure or problem associated
10 with the shared storage node can cause the SAN to go off-line and also to potentially lose
11 data that has been stored in the shared storage node. Accordingly, the basic SAN
12 configuration does not provide a high degree of data integrity and may not be acceptable
13 for use in organizations in which the risk of data loss is not acceptable.

1 SUMMARY OF THE INVENTION

2 The present invention relates to computer networks that provide virtual storage area
3 networks without using a physical shared storage node. According to the invention, the
4 network includes two or more servers, each having its own disk for storing network data.
5 In the following discussion, a network having two servers is considered. However, the
6 principles described in reference to two servers can be extrapolated to networks having
7 more than two servers.

8 When a user workstation in the network issues a write operation request to one of
9 the servers, the server receiving the request executes the write operation at its disk and uses
10 a mirror engine and a dedicated link to transmit the write operation request to other server.
11 Upon receiving the mirrored write operation request, the other server executes the write
12 operation at its disk. In this manner, data written to the disk of one server is also written to
13 the disk of another server, thereby causing the network data to be mirrored and stored at
14 both disks.

15 Since the same network data exists on the disk of both servers, either server can
16 respond to read operation requests from any user workstation. Policing protocol modules
17 at each server cooperate to implement a policing protocol, which regulates the timing and
18 priority by which each server accesses the network data. For instance, the policing
19 protocol can specify that only one server at a time can execute write requests on particular
20 portions of the network data, thereby preventing the data from being corrupted.

21 Because the data is mirrored and stored at the disk of each server in the network,
22 the network can easily tolerate the failure of one of the servers. For instance, if the first
23 server experiences a failure, the other server has access to all network data stored at its disk
24 and it can service all operation requests using its own disk. Because the same network data

1 is stored at the disk of each server in the network, the data appears, from the standpoint of
2 the servers, to have been stored in a shared storage node. Therefore, the invention provides
3 a virtual storage area network that responds operation requests and the failure of network
4 servers in a manner similar to the way in which actual storage area networks would
5 respond to failure, in that each server has immediate access to all network data.

6 The virtual storage area network and virtual shared storage nodes of the invention
7 have significant advantages compared with conventional storage area networks. For
8 instance, the networks of the invention do not require a physical shared storage node.
9 Accordingly, much of the cost associated with conventional storage area networks are
10 eliminated. Reduced costs of operating the networks of the invention make them
11 compatible with enterprises having networks with as few as two servers.

12 In addition, mirroring and storing the same network data in the disks of multiple
13 servers, in contrast to using a physical shared storage node, results in the networks of the
14 invention being significantly more tolerant of disk failure than conventional storage area
15 networks. For instance, if the disk of one of the servers of a network operated according to
16 invention were to fail, the disk of the other server in the network would have stored
17 thereon all network data. In contrast, if the physical shared storage node of a conventional
18 storage area network were to fail, the data stored thereon could be lost or, at the very least,
19 the data would be temporarily inaccessible.

20 Additional features and advantages of the invention will be set forth in the
21 description which follows, and in part will be obvious from the description, or may be
22 learned by the practice of the invention. The features and advantages of the invention may
23 be realized and obtained by means of the instruments and combinations particularly
24 pointed out in the appended claims. These and other features of the present invention will

1 become more fully apparent from the following description and appended claims, or may
2 be learned by the practice of the invention as set forth hereinafter.
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

BRIEF DESCRIPTION OF THE DRAWINGS

In order to describe the manner in which the above-recited and other advantages and features of the invention can be obtained, a more particular description of the invention briefly described above will be rendered by reference to specific embodiments thereof which are illustrated in the appended drawings. Understanding that these drawings depict only typical embodiments of the invention and are not therefore to be considered to be limiting of its scope, the invention will be described and explained with additional specificity and detail through the use of the accompanying drawings in which:

Figure 1 illustrates a network system that includes a conventional storage area network having a physical shared storage node.

Figure 2 illustrates an exemplary system that provides a suitable operating environment for the present invention.

Figure 3 is a schematic diagram illustrating a fault-tolerant network according to the invention, and shows a virtual shared storage node.

Figure 4 is a schematic diagram illustrating the fault-tolerant network of Figure 3, showing the hardware and other components that provides the functionality of the virtual shared storage node of Figure 3.

Figure 5 is a schematic diagram depicting a network according to the invention having three servers.

Figures 6 and 7 illustrate methods for mirroring network data between disks associated with two servers, thereby providing each of the servers with access to the network data.

1 Figure 8 is a flow diagram depicting a method for mirroring data between disks
2 associated with two servers, thereby providing each of the servers with access to the
3 network data.

4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

1 **DETAILED DESCRIPTION OF THE INVENTION**

2 The present invention relates to networks in which network data is mirrored and
3 stored on disks of multiple servers, such that the multiple servers provide a virtual storage
4 area network without having a physical shared storage node. Each of the multiple servers
5 in the network has a disk on which network data is stored and a mirror engine enabling the
6 server to communicate with other servers in the network. When the server receives a write
7 operation request, the server executes the write operation at its disk and transmits the write
8 operation request to the other servers in the network using the mirror engine and the
9 dedicated link or other means for communicating. The other servers receive the write
10 operation request and execute the write operation at the disks of their corresponding
11 servers. In this way, the same network data is stored at the disks of each of the multiple
12 servers. In the case of failure of one of the servers or the disk associated with any server,
13 the other server or servers remaining in the network can provide network services for any
14 user workstation in the network using the network data stored in the disks corresponding to
15 such servers.

16
17 **A. Exemplary Operating Environments**

18 The embodiments of the present invention may comprise a special purpose or
19 general-purpose computer including various computer hardware, as discussed in greater
20 detail below. Embodiments within the scope of the present invention also include
21 computer-readable media for carrying or having computer-executable instructions or data
22 structures stored thereon. Such computer-readable media can be any available media that
23 can be accessed by a general purpose or special purpose computer. By way of example,
24 and not limitation, such computer-readable media can comprise RAM, ROM, EEPROM,

1 CD-ROM or other optical disk storage, magnetic disk storage or other magnetic storage
2 devices, or any other medium which can be used to carry or store desired program code
3 means in the form of computer-executable instructions or data structures and which can be
4 accessed by a general purpose or special purpose computer. When information is
5 transferred or provided over a network or another communications connection (either
6 hardwired, wireless, or a combination of hardwired or wireless) to a computer, the
7 computer properly views the connection as a computer-readable medium. Thus, any such
8 connection is properly termed a computer-readable medium. Combinations of the above
9 should also be included within the scope of computer-readable media. Computer-
10 executable instructions comprise, for example, instructions and data which cause a general
11 purpose computer, special purpose computer, or special purpose processing device to
12 perform a certain function or group of functions.

13 Figure 2 and the following discussion are intended to provide a brief, general
14 description of a suitable computing environment in which the invention may be
15 implemented. Although not required, the invention will be described in the general context
16 of computer-executable instructions, such as program modules, being executed by
17 computers in network environments. Generally, program modules include routines,
18 programs, objects, components, data structures, etc. that perform particular tasks or
19 implement particular abstract data types. Computer-executable instructions, associated
20 data structures, and program modules represent examples of the program code means for
21 executing steps of the methods disclosed herein. The particular sequence of such
22 executable instructions or associated data structures represents examples of corresponding
23 acts for implementing the functions described in such steps.

24

1 Those skilled in the art will appreciate that the invention may be practiced in
2 network computing environments with many types of computer system configurations,
3 including personal computers, hand-held devices, multi-processor systems,
4 microprocessor-based or programmable consumer electronics, network PCs,
5 minicomputers, mainframe computers, and the like. The invention may also be practiced
6 in distributed computing environments where tasks are performed by local and remote
7 processing devices that are linked (either by hardwired links, wireless links, or by a
8 combination of hardwired or wireless links) through a communications network. In a
9 distributed computing environment, program modules may be located in both local and
10 remote memory storage devices.

11 With reference to Figure 2, an exemplary system for implementing the invention
12 includes a general purpose computing device in the form of a conventional computer 20,
13 including a processing unit 21, a system memory 22, and a system bus 23 that couples
14 various system components including the system memory 22 to the processing unit 21.
15 The system bus 23 may be any of several types of bus structures including a memory bus
16 or memory controller, a peripheral bus, and a local bus using any of a variety of bus
17 architectures. The system memory includes read only memory (ROM) 24 and random
18 access memory (RAM) 25. A basic input/output system (BIOS) 26, containing the basic
19 routines that help transfer information between elements within the computer 20, such as
20 during start-up, may be stored in ROM 24.

21 The computer 20 may also include a magnetic hard disk drive 27 for reading from
22 and writing to a magnetic hard disk 39, a magnetic disk drive 28 for reading from or
23 writing to a removable magnetic disk 29, and an optical disk drive 30 for reading from or
24 writing to removable optical disk 31 such as a CD-ROM or other optical media. Any of

1 the foregoing structures represent examples of storage devices or storage volumes that can
2 be used to establish virtual storage area networks as described herein. The magnetic hard
3 disk drive 27, magnetic disk drive 28, and optical disk drive 30 are connected to the system
4 bus 23 by a hard disk drive interface 32, a magnetic disk drive-interface 33, and an optical
5 drive interface 34, respectively. The drives and their associated computer-readable media
6 provide nonvolatile storage of computer-executable instructions, data structures, program
7 modules and other data for the computer 20. Although the exemplary environment
8 described herein employs a magnetic hard disk 39, a removable magnetic disk 29 and a
9 removable optical disk 31, other types of computer readable media for storing data can be
10 used, including magnetic cassettes, flash memory cards, digital versatile disks, Bernoulli
11 cartridges, RAMs, ROMs, and the like.

12 Program code means comprising one or more program modules may be stored on
13 the hard disk 39, magnetic disk 29, optical disk 31, ROM 24 or RAM 25, including an
14 operating system 35, one or more application programs 36, other program modules 37, and
15 program data 38. A user may enter commands and information into the computer 20
16 through keyboard 40, pointing device 42, or other input devices (not shown), such as a
17 microphone, joy stick, game pad, satellite dish, scanner, or the like. These and other input
18 devices are often connected to the processing unit 21 through a serial port interface 46
19 coupled to system bus 23. Alternatively, the input devices may be connected by other
20 interfaces, such as a parallel port, a game port or a universal serial bus (USB). A monitor
21 47 or another display device is also connected to system bus 23 via an interface, such as
22 video adapter 48. In addition to the monitor, personal computers typically include other
23 peripheral output devices (not shown), such as speakers and printers.

24

1 The computer 20 may operate in a networked environment using logical
2 connections to one or more remote computers, such as remote computers 49a and 49b.
3 Remote computers 49a and 49b may each be another personal computer, a server, a router,
4 a network PC, a peer device or other common network node, and typically include many or
5 all of the elements described above relative to the computer 20, although only memory
6 storage devices 50a and 50b and their associated application programs 36a and 36b have
7 been illustrated in Figure 2. The logical connections depicted in Figure 2 include a local
8 area network (LAN) 51 and a wide area network (WAN) 52 that are presented here by way
9 of example and not limitation. Such networking environments are commonplace in office-
10 wide or enterprise-wide computer networks, intranets and the Internet.

11 When used in a LAN networking environment, the computer 20 is connected to the
12 local network 51 through a network interface or adapter 53. When used in a WAN
13 networking environment, the computer 20 may include a modem 54, a wireless link, or
14 other means for establishing communications over the wide area network 52, such as the
15 Internet. The modem 54, which may be internal or external, is connected to the system bus
16 23 via the serial port interface 46. In a networked environment, program modules depicted
17 relative to the computer 20, or portions thereof, may be stored in the remote memory
18 storage device. It will be appreciated that the network connections shown are exemplary
19 and other means of establishing communications over wide area network 52 may be used.

20

21 B. Virtual Storage Area Networks

22 Figure 3 illustrates a representative configuration of a virtual storage area network
23 of the invention. For purposes of illustration, two server computers 310 and 320 provide
24 network services for network 301 in the example of Figure 3. However, the architecture of

1 the networks of the invention can be readily scaled to networks having three or more
2 servers, examples of which are discussed below in reference to Figure 5. Network 301
3 also includes any number of user workstations 302, which can be personal computers or
4 any other computing device that receives network services from servers 310 and 320.

5 Each server 310 and 320 includes a policing protocol module 311, 321 and an
6 input/output device driver 313, 323. Server A 310 and server B 320 operate together to
7 establish a virtual shared storage node 340. Virtual shared storage node 340 is not a
8 physical shared storage node, such as shared storage node 140 of Figure 1. Instead, virtual
9 shared storage node 340 includes various hardware and software components that will be
10 described in greater detail in reference to Figure 4, which provide functionality that, from
11 the standpoint of the I/O drivers 313 and 323, appear to be associated with an actual shared
12 storage node. Thus, in this sense, network 301 includes a virtual shared storage node 340
13 and the portion of network 301 that enables servers 310 and 320 to access virtual shared
14 storage node 340 represents a virtual storage area network.

15 It is also noted that the network configuration, including the hardware and
16 software, of network 301 outside of the region of Figure 3 designated as virtual shared
17 storage node 340 can be similar to or substantially identical to the corresponding
18 components of existing networks and similar to or substantially identical to the
19 corresponding components of conventional actual storage area networks, such as the
20 network 101 illustrated in Figure 1. One advantage of the networks operated according to
21 the invention is that they are compatible with existing policing protocol software and other
22 software that currently exists for use with conventional storage area networks.

23 Referring to Figure 4, it can be seen that the physical components of the virtual
24 shared storage node 340 are significantly different from the components of actual shared

storage node 140 of the conventional storage area network of Figure 1. For instance, network 301 of Figure 4 does not include a physical shared storage node, which eliminates much of the expense associated with acquiring and operating a storage area network. Instead, server 310 has its own disk 319, while server 320 has its own disk 329. Thus, servers 310 and 320 can be ordinary or conventional network servers, each having its own disk. In other words, the hardware of servers 310 and 320 can be similar or identical to the hardware of the majority of servers that are currently used in enterprise networks other than actual storage area networks. As used herein, "disk" and "mass storage device" are to be interpreted broadly to include any device or structure for storing data to perform the methods described herein.

The components that enable such servers 310 and 320 to provide the functionality of a virtual storage area network include mirror engines 317 and 327 and dedicated link 315. Mirror engines 317 and 327 represent examples of means for mirroring data between mass storage devices or disks of different servers. Other structures that correspond to means for mirroring data can also be used with the invention to perform the functions described herein. Moreover, as noted above, policing protocol modules 311 and 321 and other software operating at servers 310 and 320 outside of the region of Figure 4 designated as virtual shared storage node 340 can be similar or identical to existing software that has been conventionally used to operate policing protocols and other functions of actual storage area networks.

Figures 6 and 7 illustrate the manner in which mirror engines 317 and 327 and dedicated link 315 can be used to mirror data on disks 319 and 329, thereby enabling each of servers 310 and 320 to have access to all network data. As shown in Figure 6, a user of user workstation 302a causes the user workstation to issue a write operation request for

1 writing a data block A 350 to a disk associated with network 301. As shown in Figure 6,
2 the write operation request is transmitted through network 301 to server A 310. Since any
3 of the servers (e.g. server A 310 or server B 320) of network 301 can process all operation
4 requests from any user workstation, the manner in which a particular server 310 or 320 is
5 selected to process this operation is not critical to the invention. In order to balance the
6 load between servers 310 and 320, any desired load balancing algorithms can be
7 implemented. Alternatively, particular servers can be assigned to particular user
8 workstations on a default basis.

9 In this example, server A 310 receives the write operation request and passes the
10 request to I/O driver 313. I/O driver 313 then transmits the write operation request to what
11 could be perceived, from the standpoint of I/O driver 313, as a virtual shared storage node
12 (i.e., virtual shared storage node 340 of Figure 4).

13 Policing protocol module 311 operates with policing protocol module 321 of server
14 B 320 to determine whether server A 310 currently enjoys write access to disks 319 and
15 329. One primary purpose of policing protocol modules 311 and 321 is to ensure that no
16 more than a single server has write access to particular sectors or data block in disks 319
17 and 329 at any single time. Since each server 310, 320 typically has access to all network
18 data, allowing any server to have write access to the disks at all times without
19 implementing the policing protocols could otherwise lead to data corruption. Because,
20 from the standpoint of I/O drivers 313 and 323 and policing protocol modules 311 and
21 321, server A 310 and server B 320 appear to use a virtual shared storage node, the
22 policing protocols used with the invention can be similar or identical to policing protocols
23 conventionally used with actual storage area networks. In other words, as has been
24 previously mentioned, much of the software operating on server A 310 and server B 320

1 can be similar or identical to the corresponding software used with actual storage area
2 networks.

3 Since conventional policing protocols can be used with the invention, the nature of
4 policing protocols will be understood by those skilled in the art. In general, policing
5 protocols, whether used with conventional storage area networks or the virtual storage area
6 networks of the invention, determine whether a server having received an I/O request
7 currently has access priority with respect to the other servers in the network. For instance,
8 if server A 310 were to receive a write operation request, servers A 310 and servers B 320
9 communicate one with another over the network infrastructure of network 301 and use
10 policing protocol modules 311 and 321 to determine which server has write access priority
11 to the sector or other portion of the disks that is to receive the write operation. While many
12 types of policing protocols can be used with the invention, many policing protocols have
13 the common feature that they are distributed between the multiple servers and are executed
14 as the multiple servers communicate one with another.

15 Returning now to Figure 6, I/O driver transmits the write operation request to what
16 is perceived by the I/O driver as being a virtual shared storage node. Physically, however,
17 the write operation request is received by mirror engine 317. The write operation request
18 is transmitted from mirror engine 317 to disk 319, where it is executed, resulting in data A
19 350 being written to a particular sector or other region of the disk. In order to mirror data
20 A to disk 329, mirror engine 317 also transmits the write operation request to a
21 corresponding mirror engine 327 of server B 320. The write operation request is
22 transmitted by dedicated link 315 or another means for communicating. Other means for
23 communicating may include one or combination of any wireless or hard-wire
24 communication means. If dedicated link 315 is used, the physical dedicated link represents

1 one physical difference between network 301 and conventional networks. Other
2 embodiments of the invention are capable of mirroring the data using mirror engines 317
3 and 327 without dedicated link 315. For instance, the network infrastructure of network
4 301 that is otherwise used to transmit data between user workstations 302 and servers 310
5 and 320 can be used also to transmit mirrored write operation requests from one server to
6 another. Dedicated link 315, the network infrastructure of network 301 and other means
7 for communicating represent examples of means for communicating between servers and
8 the mass storage devices or disks thereof.

9 In any event, mirror engine 327 receives the mirrored write operation request and
10 transmits it to disk 329, where it is executed, resulting in data A 350 being written to disk
11 329. In this manner, after user workstation 302a issues the write operation request, the
12 data associated with the write operation request is written to disk 319 and disk 329, such
13 that both disks include mirrored copies of the same network data. It is also noted that a
14 similar process is performed when one of the user workstations 302a-n issues a write
15 operation request that causes data to be deleted from a file or otherwise deleted from disk
16 319. In other words, if data is deleted from disk 319, the same data is deleted from disk
17 329, such that the same network data is mirrored and stored at both disks.

18 Figure 7 illustrates another instance of data being mirrored between disks 319 and
19 329. In this case, user workstation 302n transmits a write operation request for writing
20 data B 352 to disk 329. The write operation request is transmitted to server B 320, where
21 it is processed by I/O driver 323 and policing protocol module 321 in a manner similar to
22 that described above in reference to the write operation request for data A 350 of Figure 6.
23 In Figure 7, mirror engine 327 transmits the write operation request to disk 329, where it is
24 executed, resulting in data B being written to disk 329. In addition, mirror engine 327

1 transmits the write operation requests to the corresponding mirror engine 317 of server A
2 310 by dedicated link 315 or by another means for communicating associated with network
3 301. Mirror engine 317 then transmits the mirrored write operation request to disk 319,
4 where it is executed, resulting in data B 352 being written to disk 319. Thus, Figures 6 and
5 7 illustrate how any server 310, 320 in network 301 is capable of causing data associated
6 with write operation requests to be stored and mirrored at each of the disks 319, 329 of the
7 servers in the network.

8 Figure 8 summarizes the method whereby write operation requests are processed
9 and the associated data is mirrored to each of the disks in the network. In step 710 the user
10 workstation issues a write operation request. In step 720, the request is transmitted over
11 the network to a particular server. In step 730, the particular server executes the policing
12 protocols to determine whether the server currently has write access to the disks in the
13 network. If, according to decision block 732, the server does not have write access, the
14 operation request is denied in step 751 and aborted in step 753. Alternatively, if, according
15 to decision block 732, the server does not have write access, the operation can be queued
16 until such time that write access is granted to the server.

17 If, according to decision block 732, the server does have write access, the operation
18 request is accepted at step 741. The mirror engine then copies the write operation request
19 in step 743 and transmits it to one or more other servers in the network. The particular
20 server that received the write operation request executes the write operation at its disk in
21 step 749. At the same time or just prior to or after step 749, step 750 is executed, in which
22 the other server or servers in the network, which have received the mirrored copy of the
23 write operation request, execute the mirrored write operation request, such that the same

24

1 network data is also stored in the disks associated with the other server or servers. The
2 order in which steps 743, 749 and 750 are conducted is not critical to the invention.

3 The manner in which network 301 can respond to and tolerate the failure of a
4 server or disk is illustrated in reference to Figure 7. In this example, it is assumed that data
5 A 350 and data B 352 have already been mirrored to both disks 319 and 329 as described
6 above. Prior to server A 310 going offline, server A advertises the fact that it can provide
7 access for workstations to the virtual shared storage node, illustrated at 340 of Figure 4,
8 such that any workstation requiring services of the disks of the virtual shared storage node
9 can receive them through server A 310. In normal operation, server A 310 provides access
10 to the virtual shared storage node by using disk 319. In this example, it is assumed that
11 server A 310 remains operational, but its associated disk 319 goes offline and becomes
12 unavailable for any reason. Server A 310 continues to advertise to workstations that it can
13 provide access to the virtual shared storage node, including disk 319 because, from the
14 standpoint of server A 310, the virtual shared storage node and the data stored thereon
15 remain available.

16 After the failure of disk 319, workstations 302a-d can continue to issue read
17 operation requests to be processed by the virtual shared storage node through server A 310.
18 In this example, it is assumed that workstation 302a issues a read operation request
19 directed to data A 350. Upon the read operation request being received by server A 310,
20 the read operation request is received by I/O driver 313 and transmitted to mirror engine
21 317, which, as shown at Figure 4, is within virtual shared storage node 340.

22 At this point, the read operation request has transmitted in the typical manner to a
23 storage device that is perceived, from the standpoint of server A 310, as being a shared
24 storage node. However, as mentioned above, disk 319 is not accessible and cannot service

1 the read operation request. Accordingly, the read operation request is transmitted to server
2 B 320 using dedicated link 315. The read operation request is then used to access disk
3 329, which has a full copy of the network data, including data A 350. Thus, network 301
4 is capable of seamlessly responding to inaccessibility of disk 319 by using mirror engines
5 317 and 327 to redirect read operation requests that are received by server A 310.
6 Operation of network 301 continues uninterrupted notwithstanding the failure of disk 319.
7 Moreover, server A 310 can respond to other network operation requests, such as write
8 operation requests, in a similar manner after the failure of disk 319 by using the virtual
9 shared storage node.

10 The foregoing method of responding to disk failure enables network activity to
11 continue without disruption of any network activity that could have been partially
12 completed at the time of the disk failure. Responding to disk failure in this way requires
13 an operational I/O driver 313 and mirror engine 317 of server A 310.

14 If these functional components of server A 310 become inoperable, network 301
15 has a secondary way of continuing to provide access to network data according to one
16 embodiment. In this scenario, if user workstation 302a were to issue a read operation
17 request that would otherwise be processed by server A 310, the read operation request can
18 be serviced by server B 320, since server B 320 has access to all network data on its disk
19 329. For purposes of illustration, it is assumed that the read operation request issued by
20 user workstation 302 is directed to data A 350. Because server A 310 is offline, server B
21 320 processes the read operation request. Server B 320 uses the mirrored copy of the
22 network data stored at disk 329 to service the read operation request and thereby provide
23 user workstation with read access to data A 350. It is noted that conventional storage area
24 networks also enable all servers to provide read access to all network data in the case of

1 one of the servers of the network experiencing a failure or otherwise going offline.
2 However, unlike conventional storage area networks, the networks of the invention do not
3 use a physical shared storage node to provide access to all network data through any
4 server.

5

6

7 The foregoing examples of the capability of the networks of the invention to
8 continue operating after disk or server failure provide significant advantages that are not
9 possible using a conventional storage area network. Typically, a conventional storage area
10 network has a single component that, if it fails, can render the data inaccessible. For
11 instance a typical conventional storage area network includes a SAN connection card or a
12 disk driver that must be operational in order to provide access to the shared storage node.

13 In addition, physical failure of the disks of the shared storage node of a
14 conventional storage area network can cause the loss of access to the data. Indeed, if
15 shared storage node 140 of Figure 1 were to be physically damaged or if data stored
16 thereon were to be physically lost, the conventional storage area network might experience
17 permanent and irretrievable loss of network data in addition to the down time associated
18 with the failed shared storage node 140. In order to eliminate this possibility,
19 administrators of conventional storage area networks can purchase and maintain in the
20 shared storage node redundant arrays of independent disks (RAIDs), which increases the
21 cost and complexity of the system. As described above, the present invention enables fault
22 tolerance of disks in a virtual shared storage node without requiring RAIDs.

23

24

1 The methods of invention illustrated in Figures 3, 4, 6 and 7 in reference to two
2 servers can be scaled to networks having more than two servers. For instance, Figure 5
3 illustrates a network according to the invention having three servers, namely, server A 520,
4 server B 520 and server C 530. The operation of network 501 of Figure 5 is similar to that
5 of network 301 described above in reference to Figures 6 and 7. For instance, when server
6 A receives a write operation request from a user workstation 502, the write operation
7 request is processed by I/O driver 513. Policing protocol module 511 and server A 510
8 operate in combination with policing protocol module 521 of server B 520 and policing
9 protocol module 531 of server C 530.

10 A mirror engine 517 of server A 510 transmits a copy of the write operation request
11 to mirror engine 527 of server B 520 through dedicated link 515 or other another
12 communications link. Mirror engine 517 also transmits a copy of the write operation
13 request to mirror engine 537 of server C 530 through dedicated link 555 or other
14 communications link. Again, it is noted that any other communications link can be used to
15 transmit the copies of the write operation request to the various mirror engines of the other
16 servers in the network. For instance, the network infrastructure of network 501 can be
17 used to transmit such write operation request. Alternatively, a write operation request can
18 be transmitted from mirror engine 517 to mirror engine 537 by transmitting the operation
19 request sequentially through dedicated link 515, mirror engine 527 and dedicated link 525.
20 All that is important is that mirror engines 517, 527, and 537 be capable of communicating
21 one with another. In the foregoing manner, data written to one of the disks 519, 529 and
22 539 is stored at all the disks. In the case of the failure of one of servers 510, 520, 530,
23 remaining servers are capable of servicing all requests from any user workstations for any
24 of the network data.

1 The present invention may be embodied in other specific forms without departing
2 from its spirit or essential characteristics. The described embodiments are to be considered
3 in all respects only as illustrative and not restrictive. The scope of the invention is,
4 therefore, indicated by the appended claims rather than by the foregoing description. All
5 changes which come within the meaning and range of equivalency of the claims are to be
6 embraced within their scope.

7 What is claimed and desired to be secured by United States Letters Patent is:

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24